



# SAM 奇怪的补充

KF.Dong

# SAM



- 只能识别一个串**所有后缀**的自动机；
- 识别的定义：从自动机的源开始，依此读入一个字符然后沿着自动机的边转移，若到达的状态是结束状态，则此串可以被识别。

# 关于定义和约定

- 如果未说明，大写 $S$ 表示原串， $\text{Fac}$ 表示 $S$ 的所有连续子串， $\text{Suf}$ 表示 $S$ 的所有后缀， $\text{Suffix}(a)$ 表示 $S$ 从 $a$ 开始的后缀， $S[l,r)$ 表示 $S$ 的一个子串
- 对于两个字符串 $a$ 、 $b$ ，定义 $ab$ 为将 $b$ 接在 $a$ 之后的新的字符串
- $\text{ST}(s)$ 表示在 $S$ 的自动机上从起点出发，读入一个字符串 $s$ 后，在自动机 $S$ 上到达的状态(state)
- 根据后缀自动机的定义，若 $\text{ST}(s) \in \text{end}$ (结束状态)，则 $s$ 为 $S$ 的一个后缀；

# 一些结论

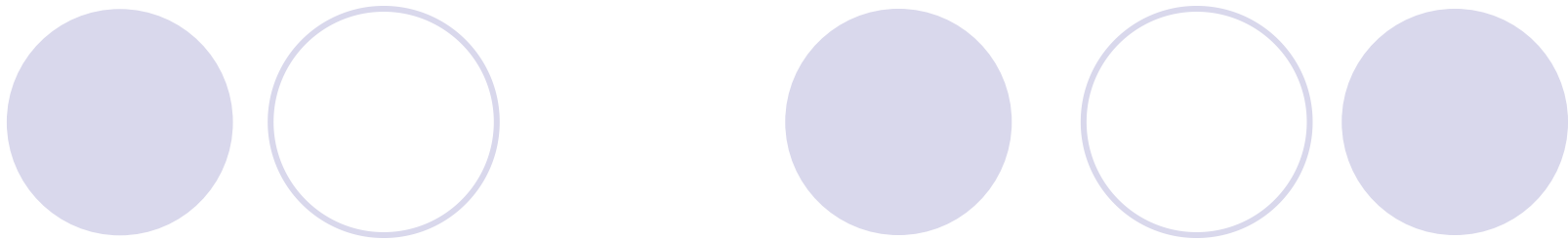
- 若 $s$ 不是 $S$ 的一个子串，那么 $ST(s)=NULL$ ；因为 $s$ 在后面加上任何字符都不可能是 $S$ 的后缀，所以我们不需要记录这个状态。
- 因此当且仅当 $ST(s) \neq NULL$ 时， $s$ 是 $S$ 的一个子串
- 考虑从 $ST(a)$ 开始，能识别哪些串，即有哪些 $x$ ，使得 $ax \in Suffix$ 。显然 $x \in Suffix$ ，故 $ST(a)$ 能识别的串还是一些后缀，记 $Reg(ST(a)) = \{x \mid ax \in Suffix\}$

# 分析

- 由于 $\text{Reg}()$ 是一些后缀的集合，即 $\text{Reg}(\text{ST}(a)) = \{\text{Suffix}(r_1), \text{Suffix}(r_2), \dots, \text{Suffix}(r_n)\}$
- 不妨记 $\text{Right}'(a) = \{r_1, r_2, \dots, r_n\}$ ，也就是 $\text{Right}'(a)$ 和 $\text{Reg}(\text{ST}(a))$ 是对应的，可以看出 $\text{Right}'(a)$ 就是 $a$ 在 $S$ 中出现的所有结束位置的集合
- 那么若 $\text{Right}'(a) = \text{Right}'(b)$ ，则 $\text{ST}(a) = \text{ST}(b)$ ，这两个状态的所有后继状态完全一样；
- 所以在最简状态SAM中，每一个状态 $s$ 满足：所有 $\text{Right}'(a)$ 相同的串的转移 $\text{ST}(a) = s$ ，以及所有 $\text{ST}(a) = s$ 的串 $a$ ， $\text{Right}'(a)$ 相同，不妨记 $\text{Right}(s) = \text{Right}'(a)$

# 分析

- 接下来，考虑一个状态 $s=ST(a)$ 中的字符串 $a$ 在原串 $S$ 中的分布；
- 首先，由于 $Right'(a)$ 是 $axi \in Suffix$ 的 $xi$ 的起始位置 $ri$ 的集合，例：  
●  $aaabbaaabd$  中  
 $Right'(aaab)=\{4,9\}$ ，对于一个状态 $s$ ，设 $Right(s)=\{r1...rn\}$ ，那么如果知道了 $a$ 串的长度 $len$ ，就可以确定 $a$ 在原串 $S$ 中是 $[ri-len, ri)$  ( $1 \leq i \leq n$ )，容易看出对于一个 $len$ ，串 $a$ 有且只有一个，但是有些 $len$ 是不合法的，比如 $s=ST(aaab), len=1$ 时， $a$ 串是“b”，但是 $Right'(a)=\{3,4,9\}$ ，那么 $ST(a) \neq s$
- 显然，对于一个状态 $s$ ，若长度 $l \leq r$ 所对应的 $a_l, a_r$ 满足 $ST(a_l)=s, ST(a_r)=s$ ，那么对于任意 $l \leq x \leq r, ST(a_x)=s$ ，不妨设 $s$ 所有合法的长度集合为 $[Min(s), Max(s)]$



- 由于 $\text{Right}'(a)$ 是串 $a$ 出现的结束位置的集合，那么显然对于两个串 $a, b$ ,  $\text{Right}'(a)$ 和 $\text{Right}'(b)$ 要么不相交，要么是包含关系；
- 所以可以看出 $\text{Right}()$ 是树形关系，那么记 $\text{Parent}(s)$ 是状态 $s$ 在树上的祖先。  
即 $\text{Right}(s) \in \text{Right}(\text{Parent}(s))$

举个例子：aaabbbaabd

$\text{Parent}(\text{ST}(\text{aaab})) = \text{Parent}(\text{ST}(\text{ab})) = \text{ST}(\text{b})$

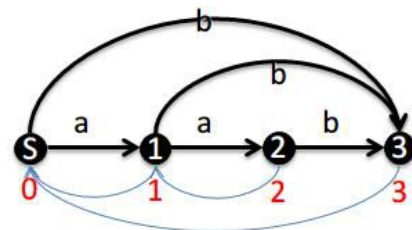
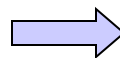
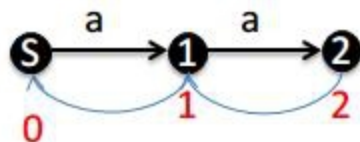
容易看出， $\text{Max}(\text{Parent}(s)) = \text{Min}(s) - 1$ ，也就是说，状态 $s$ 的长度 $\text{Min}(s) - 1$ 对应的串的出现次数比 $\text{Min}(s)$ 多了一些

由此可以证明状态数是线性的，具体证明略

另外，由于包含关系可知，所有的end就是 $\text{ST}(S)$ 的所有Parent树的祖先

# 构造算法

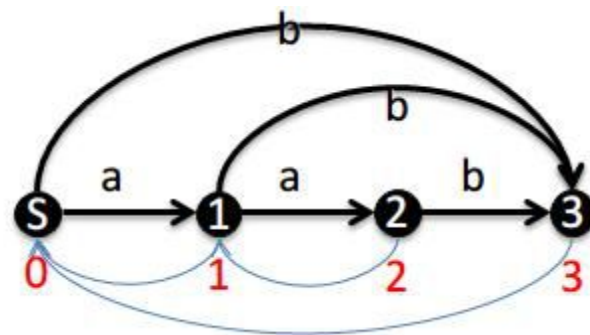
- 每次在S末尾添加一个字符x, 然后由S的SAM改造成Sx的SAM
- 设S的长度为len, 新建状态np,  $\text{Right}(np) = \text{len} + 1$ ;
- 若状态s,  $\text{Right}(s)$ 不包含len, 那么添加一个字符后, 状态s不需要改变。
- 设 $p = \text{ST}(S)$ ; 由于 $\text{Right}(p) = \{\text{len}\}$ , 那么可能改变的状态就是p在Parent树上的所有祖先, 不妨记为 $v_1 = p, v_2, v_3 \dots v_k = \text{root}$  (按深度顺序), 那么 $v_1$ 到 $v_k$ 的Right集合是递增的;
- 若 $\text{ST}(v_i, x) = \text{NULL}$ , 也就是说 $v_i$ 的结束位置之后不存在x, 那么由于 $v_i$ 代表的一个串是S的后缀 (因为 $\text{Right}(v_i)$ 包含len), 所以令 $\text{ST}(v_i, x) = np$ , 代表在新的自动机上, 若一个串转移到了状态 $v_i$ , 再读入一个字符x, 转移到的新状态就是np
- 比如原串aa, 新加入b, 那么 $\text{ST}(2, 'b') = \text{ST}(1, 'b') = \text{ST}(s, 'b') = np$  (图中蓝色边表示Parent, 红色数字表示 $\text{Max}(s)$ )





# 构造算法

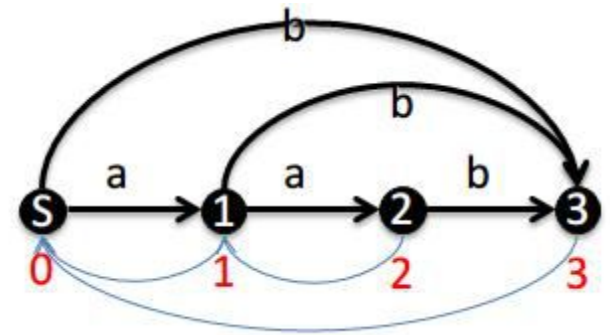
- 若  $ST(v_i, x) = q \neq \text{NULL}$ :
- 如原串是 **aab**，它的自动机是这样：



- 现在添加字符 **b**，对于状态3  
直接连一条到新状态4的**b**的边,但是对于s怎么办？

# 构造算法

- $s$  状态的 **Right** 集合中，只有那些  $S[ri]=x$  的  $ri$  仍然是合法的，但是如果直接把  $ST(s,x)$  赋为合法的 **end** 的话，在图中3就是合法的 **end** 了，然而  $ST(s, "aab")$  不是合法后缀却也赋为合法的 **end** 了（此时的串是 **aabb**）



当然，如果  $\text{Max}(vi)+1=\text{Max}(q)$ ，也就是说不存在另一个串  $t \neq 'x'$ ，使得  $ST(vi,t)=St(vi,x)$ ，那么就可以直接这样赋值为合法了。（就是令  $\text{Parent}(np)=vi$ ，即在  $\text{Right}(vi)$  中间插入  $\{len+1\}$ ）

# 构造算法

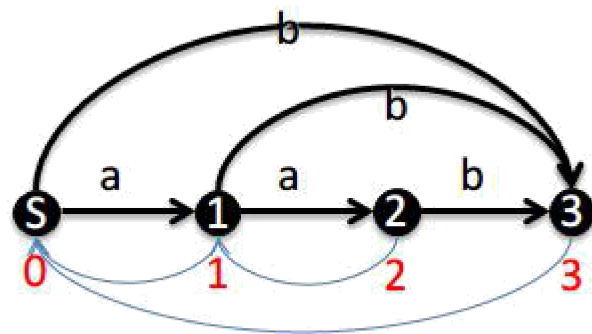
- 如果  $\text{Max}(v_i) + 1 \neq \text{Max}(q)$ :

由于不能直接赋值，那么就把  $q$  复制出一个新的节点  $\text{new}$ ，然后令  $\text{Max}(\text{new}) = \text{Max}(v_i) + 1$ ，这个时候就可以直接令  $\text{Right}(\text{new}) = \text{Right}(q) \cup \{\text{len} + 1\}$ ，就是令  $\text{Parent}(\text{np}) = \text{new}$ ；可以证明  $\text{Parent}(q) = \text{new}$ ；

随后，所有的  $\text{Right}(v_j)$  包含  $\text{len}$  的（就是  $\text{Parent}$  树上的祖先），并且  $\text{ST}(v_j, x) = q$  的  $v_j$  状态，令转移  $\text{ST}(v_j, x) = \text{new}$  即可。这是因为原来  $v_j$  可以转移到  $q$  而且  $\text{len} \in \text{Right}\{v_j\}$ ，那么在新的字符串中， $v_j$  可以转移到的状态除了  $q$  以外，还多了一个  $\{\text{len} + 1\}$

# 构造算法

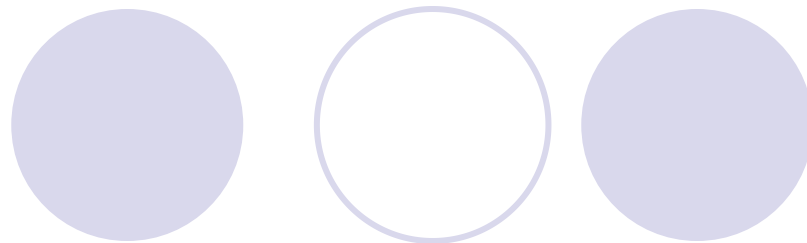
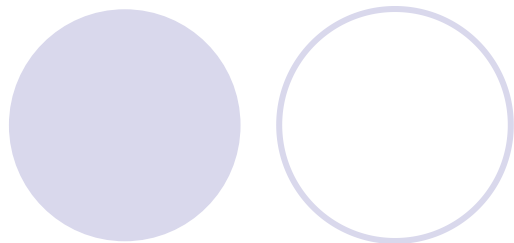
- 比如说这张图中,  $\text{Right}(\text{ST}(s,b))=\{3\}$ , 然后如果新加了一个字符**b**的话,  $\text{Right}(\text{ST}(s,b))$ 就应该变成 $\{3,4\}$ 了



- 再来一个例子: `aaabcaabaa`, 状态 $v_i$ 所代表的字符串用红色标出: `aaabcaabaa`, 这个时候加入一个字符**b**, 原来的 $q$ 代表的字符串是: `aaabcaabaa`, 现在应该是`aaabcaabaab`, 比原来多了一个 $\{len+1\}$ , 所以原来转移到 $q$ 而且 $\text{Right}$ 集合包含 $len$ 的状态 $v_j$ ,  $\text{ST}(v_j, "b")$ 就应该变成 $new$

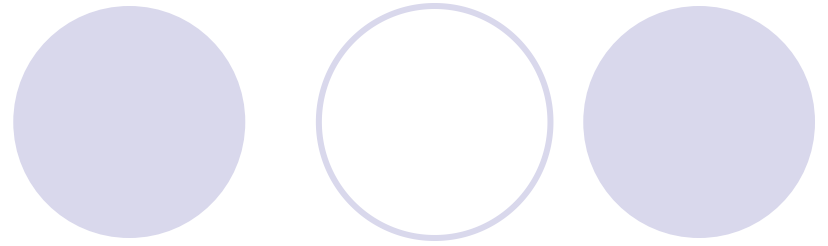
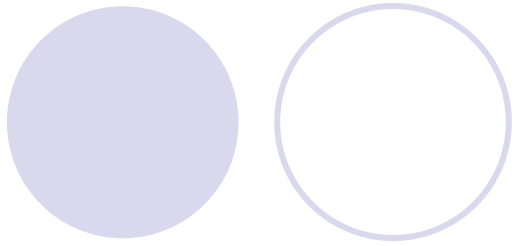
# 具体实现

- 新建节点 $np$ , 令 $Max(np)=len+1$ ;
- 设 $ST(s,S)=last$ , 那么 $p$ 从 $last$ 开始, 若 $ST(p,x)=NULL$ , 直接令 $ST(p,x)=np$ , 直到 $ST(p,x)\neq NULL$
- 若不存在这样的 $p$ , 令 $Parent(np)=s$ ; 结束
- 否则, 设 $ST(p,x)=q$ , 若 $Max(q)=Max(p)+1$ , 令 $Parent(np)=q$ ; 结束
- 新建节点 $new=q$ , 令 $Parent(q)=Parent(np)=new$ ,  $Max(new)=Max(p)+1$ , 然后再从 $p$ 开始沿着 $Parent$ 边, 若 $ST(p,x)=q$ , 就令 $ST(p,x)=q$ ;



## ● 代码并不长.....

```
void Extend(int C, int Len) {
    int Now = ++Total, P;
    SAM[Now].L = Len;
    for (P = Last; P && !SAM[P].Ch[C]; P = SAM[P].Fail)
        SAM[P].Ch[C] = Now;
    Last = Now;
    if (!P) SAM[Now].Fail = 1;
    else {
        if (SAM[SAM[P].Ch[C]].L == SAM[P].L + 1) SAM[Now].Fail = SAM[P].Ch[C];
        else {
            int New = ++Total, Q = SAM[P].Ch[C];
            SAM[New] = SAM[Q];
            SAM[New].L = SAM[P].L + 1;
            SAM[Q].Fail = SAM[Now].Fail = New;
            for (; P && (SAM[P].Ch[C] == Q); P = SAM[P].Fail)
                SAM[P].Ch[C] = New;
        }
    }
}
```



● 推荐题目：

○ BZOJ2555

○ BZOJ2806

○ BZOJ3172

○ HDU4416